In []:

PROBLEM STATEMENT: TO PREDICT AND ANALYZE WHICH GENDER

In [1]:

```
import numpy as np
import pandas as pd

from sklearn import preprocessing
import matplotlib.pyplot as plt
#plt.rc("font", size=14)
import seaborn as sns
sns.set(style="white")
sns.set(style="white")
import warnings
warnings.simplefilter(action='ignore')
```

In [2]:

train_df = pd.read_csv(r"C:\Users\smb06\OneDrive\Desktop\train.gender.csv")
train_df

Out[2]:

| | Passengerld | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fa |
|-------|---------------|----------|--------|---|--------|------|-------|-------|---------------------|-------|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.25 |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th | female | 38.0 | 1 | 0 | PC 17599 | 71.28 |
| 2 | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.92 |
| 3 | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.10 |
| 4 | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.05 |
| | | | | | | | | | | |
| 886 | 887 | 0 | 2 | Montvila, Rev. Juozas | male | 27.0 | 0 | 0 | 211536 | 13.00 |
| 887 | 888 | 1 | 1 | Graham, Miss. Margaret Edith | female | 19.0 | 0 | 0 | 112053 | 30.00 |
| 888 | 889 | 0 | 3 | Johnston, Miss. Catherine Helen "Carrie" | female | NaN | 1 | 2 | W./C. 6607 | 23.45 |
| 889 | 890 | 1 | 1 | Behr, Mr. Karl Howell | male | 26.0 | 0 | 0 | 111369 | 30.00 |
| 890 | 891 | 0 | 3 | Dooley, Mr. Patrick | male | 32.0 | 0 | 0 | 370376 | 7.75 |
| 891 r | ows × 12 colu | ımns | | | | | | | | |
| 4 | 12 0010 | | | | | | | | | |
| 1 | | | | | | | | | | |

In [3]:

test_df = pd.read_csv(r"C:\Users\smb06\OneDrive\Desktop\test.gender.csv")
test_df

Out[3]:

| | Passengerld | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cal |
|-----|-------------|--------|--|--------|------|-------|-------|-----------------------|----------|-----|
| 0 | 892 | 3 | Kelly, Mr. James | male | 34.5 | 0 | 0 | 330911 | 7.8292 | N |
| 1 | 893 | 3 | Wilkes, Mrs. James (Ellen Needs) | female | 47.0 | 1 | 0 | 363272 | 7.0000 | N |
| 2 | 894 | 2 | Myles, Mr. Thomas Francis | male | 62.0 | 0 | 0 | 240276 | 9.6875 | N |
| 3 | 895 | 3 | Wirz, Mr. Albert | male | 27.0 | 0 | 0 | 315154 | 8.6625 | N |
| 4 | 896 | 3 | Hirvonen, Mrs. Alexander (Helga E Lindqvist) | female | 22.0 | 1 | 1 | 3101298 | 12.2875 | N |
| | | | | | | | | | | |
| 413 | 1305 | 3 | Spector, Mr. Woolf | male | NaN | 0 | 0 | A.5. 3236 | 8.0500 | N |
| 414 | 1306 | 1 | Oliva y Ocana, Dona. Fermina | female | 39.0 | 0 | 0 | PC 17758 | 108.9000 | C1 |
| 415 | 1307 | 3 | Saether, Mr. Simon Sivertsen | male | 38.5 | 0 | 0 | SOTON/O.Q. 3101262 | 7.2500 | N |
| 416 | 1308 | 3 | Ware, Mr. Frederick | male | NaN | 0 | 0 | 359309 | 8.0500 | N |
| 417 | 1309 | 3 | Peter, Master. Michael J | male | NaN | 1 | 1 | 2668 | 22.3583 | N |

418 rows × 11 columns

4

In [4]:

train_df.shape

Out[4]:

(891, 12)

In [5]:

test_df.head()

Out[5]:

| | Passengerld | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Em |
|---|-------------|--------|--|--------|------|-------|-------|---------|---------|-------|----|
| 0 | 892 | 3 | Kelly, Mr. James | male | 34.5 | 0 | 0 | 330911 | 7.8292 | NaN | |
| 1 | 893 | 3 | Wilkes, Mrs. James (Ellen Needs) | female | 47.0 | 1 | 0 | 363272 | 7.0000 | NaN | |
| 2 | 894 | 2 | Myles, Mr. Thomas Francis | male | 62.0 | 0 | 0 | 240276 | 9.6875 | NaN | |
| 3 | 895 | 3 | Wirz, Mr. Albert | male | 27.0 | 0 | 0 | 315154 | 8.6625 | NaN | |
| 4 | 896 | 3 | Hirvonen, Mrs. Alexander (Helga E Lindqvist) | female | 22.0 | 1 | 1 | 3101298 | 12.2875 | NaN | |
| 4 | | | | | | | | | | | • |

In [6]:

test_df.shape

Out[6]:

(418, 11)

In [7]:

```
train_df.describe
```

Out[7]:

| | nd meth | od NDFrame. 1 2 3 4 5 887 888 889 890 891 | describ | e of 3 \ 1 \ 3 \ 2 \ 1 \ 3 \ 1 \ 3 \ 3 \ 3 \ 3 \ 3 \ 1 \ 3 \ 3 \ 1 \ 1 \ 3 \ 1 \ 1 \ 3 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 | | engerId | Surv | /ived | Pclass | |
|---------------------------|---------|--|--|--|----------------------------------|----------|----------------------------|--------|--------|------|
| | | | | | | N | ame | Sex | a Age | SibS |
| р 0 | | | | Braund | , Mr. C |)wen Har | ris | male | 22.0 | |
| 1 \ 1 | | s, Mrs. Joh | n Bradl | ey (Flor | ence Br | iggs Th | | female | 38.0 | |
| 1 2 | | | | Heikk | inen, M | liss. La | ina | female | 26.0 | |
| 0 3 | F | utrelle, Mr | s. Jacq | ues Heat | h (Lily | ⁄ May Pe | el) | female | 35.0 | |
| 1 4 | | | | Allen, | Mr. Wil | .liam He | nry | male | 35.0 | |
| 0 | | | | | | | | | • • • | |
| 886 | | | | Mont | vila, R | Rev. Juo | zas | male | 27.0 | |
| 0 887 | | | Gra | ham, Mis | s. Marg | garet Ed | ith | female | 19.0 | |
| 0 888 | | Johnston | , Miss. | Catheri | ne Hele | n "Carr | ie" | female | . NaN | |
| 1 889 | | | | Behr | , Mr. k | (arl How | ell | male | 26.0 | |
| 0 890 0 | | | | Do | oley, M | lr. Patr | ick | male | 32.0 | |
| 0 1 2 3 4 | Parch | A/5 PC STON/O2. 3 | Ticket 21171 17599 101282 113803 373450 211536 | 7.2500 71.2833 7.9250 53.1000 8.0500 | NaN C85 NaN C123 NaN | | d S C S S S | | | |
| 887 888 | 0 2 | | 112053 . 6607 | 30.0000 | B42 | | S S | | | |
| 889 890 | 0 0 | | 111369 370376 | 30.0000 7.7500 | C148 | | C Q | | | |
| _ | | | _ | | | | | | | |

[891 rows x 12 columns]>

In [8]:

```
train_df.info()
```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):

| # | Column | Non-Null Count | Dtype |
|----|-------------|----------------|---------|
| | | | |
| 0 | PassengerId | 891 non-null | int64 |
| 1 | Survived | 891 non-null | int64 |
| 2 | Pclass | 891 non-null | int64 |
| 3 | Name | 891 non-null | object |
| 4 | Sex | 891 non-null | object |
| 5 | Age | 714 non-null | float64 |
| 6 | SibSp | 891 non-null | int64 |
| 7 | Parch | 891 non-null | int64 |
| 8 | Ticket | 891 non-null | object |
| 9 | Fare | 891 non-null | float64 |
| 10 | Cabin | 204 non-null | object |
| 11 | Embarked | 889 non-null | object |
| | | | |

dtypes: float64(2), int64(5), object(5)

memory usage: 83.7+ KB

In [9]:

test_df.describe

Out[9]:

| | | d NDFr | ame.des | cribe of | PassengerId | Pclass | |
|----------|--------|------------|---------|----------|----------------------|-------------|--------------|
| Name | | 902 | 2 | | | Vally M | In James \ |
| 0 1 | | 892 893 | 3 3 | | Wilkes, Mrs. J | | lr. James \ |
| 2 | | 894 | 2 | | | Mr. Thomas | |
| 3 | | | 3 | | Myres, | | |
| 3 4 | | 895 | | 114 | - Mar Alayandan (| Wirz, Mr | |
| 4 | | 896 | 3 | Hirvonei | n, Mrs. Alexander (| neiga E Li | .naqvist) |
| 413 | | 1305 | 3 | | | Spector, M | lr. Woolf |
| 414 | | 1306 | 1 | | Oliva y Oc | • | |
| 415 | | 1307 | 3 | | | lr. Simon S | |
| 416 | | 1308 | 3 | | | lare, Mr. F | |
| 417 | | 1309 | 3 | | | Master. M | |
| , | | 2303 | | | 1 0001) | | izenacz 5 |
| | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin Embark |
| ed | _ | | | _ | | | |
| 0 | male | 34.5 | 0 | 0 | 330911 | 7.8292 | NaN |
| Q | | | _ | _ | | | |
| 1 | female | 47.0 | 1 | 0 | 363272 | 7.0000 | NaN |
| S | _ | | | _ | | | |
| 2 | male | 62.0 | 0 | 0 | 240276 | 9.6875 | NaN |
| Q | , | | | • | 245454 | 0.6605 | |
| 3 | male | 27.0 | 0 | 0 | 315154 | 8.6625 | NaN |
| S | | | _ | | | | |
| 4 | female | 22.0 | 1 | 1 | 3101298 | 12.2875 | NaN |
| S | | | | | | | |
| • • | • • • | • • • | • • • | • • • | ••• | ••• | • • • |
| 413 | male | NaN | 0 | 0 | A.5. 3236 | 8.0500 | NaN |
| 413 S | mare | ivaiv | Ø | V | A.5. 3230 | 0.0300 | Ivaiv |
| 3 414 | female | 20.0 | 0 | 0 | DC 177E0 | 100 0000 | C10E |
| 414 C | тешате | 39.0 | 0 | 0 | PC 17758 | 108.9000 | C105 |
| | | 20 5 | 0 | 0 (| COTON (O. O. 2404262 | 7 2500 | NI - NI |
| 415 | male | 38.5 | 0 | 0 9 | SOTON/O.Q. 3101262 | 7.2500 | NaN |
| S | , | | • | • | 250200 | 0.0500 | |
| 416 | male | NaN | 0 | 0 | 359309 | 8.0500 | NaN |
| S | - | | _ | | ** | 00 0-0- | |
| 417 | male | NaN | 1 | 1 | 2668 | 22.3583 | NaN |
| C | | | | | | | |

[418 rows x 11 columns]>

In [10]:

```
test_df.info()
```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 418 entries, 0 to 417
Data columns (total 11 columns):
 # Column Non-Null Count Dty

| # | Column | Non-Null Count | Dtype |
|------|---------------|----------------|---------|
| | | | |
| 0 | PassengerId | 418 non-null | int64 |
| 1 | Pclass | 418 non-null | int64 |
| 2 | Name | 418 non-null | object |
| 3 | Sex | 418 non-null | object |
| 4 | Age | 332 non-null | float64 |
| 5 | SibSp | 418 non-null | int64 |
| 6 | Parch | 418 non-null | int64 |
| 7 | Ticket | 418 non-null | object |
| 8 | Fare | 417 non-null | float64 |
| 9 | Cabin | 91 non-null | object |
| 10 | Embarked | 418 non-null | object |
| dtvn | es: float64(2 |) int64(4) ohi | ect(5) |

dtypes: float64(2), int64(4), object(5)

memory usage: 36.1+ KB

To find Missing values

In [11]:

train_df.isnull().sum()

Out[11]:

PassengerId 0 Survived 0 Pclass 0 Name Sex 0 Age 177 SibSp 0 Parch 0 Ticket Fare 0 Cabin 687 2 Embarked dtype: int64

In [12]:

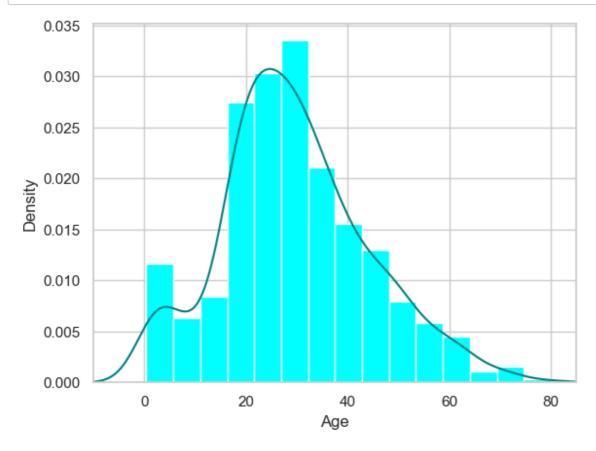
```
test_df.isnull().sum()
```

Out[12]:

PassengerId 0 **Pclass** 0 Name 0 Sex 0 86 Age SibSp 0 Parch 0 Ticket 0 Fare 1 Cabin 327 Embarked 0 dtype: int64

In [13]:

```
ax=train_df["Age"].hist(bins=15, density=True, stacked=True,color='cyan')
train_df['Age'].plot(kind='density', color='teal')
ax.set(xlabel='Age')
plt.xlim(-10,85)
plt.show()
```



```
In [14]:
```

```
print(train_df["Age"].mean(skipna=True))
print(train_df["Age"].median(skipna=True))
```

29.69911764705882

28.0

In [15]:

```
print((train_df['Cabin'].isnull().sum()/train_df.shape[0])*100)
```

77.10437710437711

In [16]:

```
print((train_df['Embarked'].isnull().sum()/train_df.shape[0])*100)
```

0.22446689113355783

In [17]:

```
print('Board passengers grouped by port of embarkation (c = cherbourg, Q = Queenstown)')
print(train_df['Embarked'].value_counts())
sns.countplot(x='Embarked', data=train_df, palette='Set2')
plot.show()
```

```
Board passengers grouped by port of embarkation (c = cherbourg, Q = Queens town)

Embarked

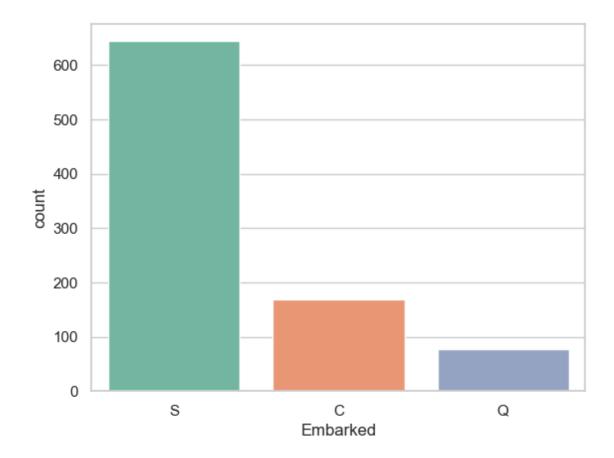
S 644

C 168

Q 77

Name: count, dtype: int64
```

NameError: name 'plot' is not defined



In [18]:

```
print(train_df['Embarked'].value_counts().idxmax())
```

S

In [19]:

```
train_data = train_df.copy()
train_data["Age"].fillna(train_df["Age"].median(skipna=True),inplace=True)
train_data["Embarked"].fillna(train_df['Embarked'].value_counts().idxmax(),inplace=True)
train_data.drop('Cabin', axis=1, inplace=True)
```

In [20]:

```
train_data.isnull().sum()
```

Out[20]:

| PassengerId | 0 |
|--------------|---|
| Survived | 0 |
| Pclass | 0 |
| Name | 0 |
| Sex | 0 |
| Age | 0 |
| SibSp | 0 |
| Parch | 0 |
| Ticket | 0 |
| Fare | 0 |
| Embarked | 0 |
| dtype: int64 | |
| | |

In [21]:

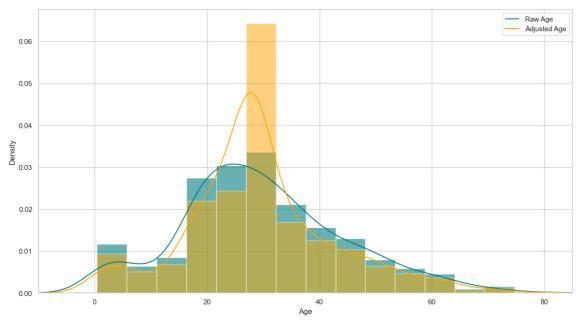
```
train_data.head()
```

Out[21]:

| | Passengerld | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare |
|---|-------------|----------|--------|---|--------|------|-------|-------|---------------------|-------------|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 |
| 2 | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 |
| 3 | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 |
| 4 | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.0500 |
| 4 | | | | | | | | | | > |

In [22]:

```
plt.figure(figsize=(15,8))
ax = train_df["Age"].hist(bins=15, density=True, stacked=True, color='teal', alpha=0.6)
train_df["Age"].plot(kind='density', color='teal')
ax = train_data["Age"].hist(bins=15, density=True, stacked=True, color='orange', alpha=0
train_data["Age"].plot(kind='density', color='orange')
ax.legend(['Raw Age', 'Adjusted Age'])
ax.set(xlabel='Age')
plt.xlim(-10,85)
plt.show()
```



In [23]:

```
train_data['TravelAlone']=np.where((train_data["SibSp"]+train_data)["Parch"]>0,0,1)
train_data.drop('SibSp',axis=1, inplace=True)
train_data.drop('Parch',axis=1, inplace=True)
```

In [24]:

```
training=pd.get_dummies(train_data, columns=["Pclass","Embarked","Sex"])
training.drop('Sex_female', axis=1, inplace=True)
training.drop('PassengerId', axis=1, inplace=True)
training.drop('Name', axis=1, inplace=True)
training.drop('Ticket', axis=1, inplace=True)

final_train = training
final_train.head()
```

Out[24]:

| | Survived | Age | Fare | TravelAlone | Pclass_1 | Pclass_2 | Pclass_3 | Embarked_C | Embark |
|---|----------|------|---------|-------------|----------|----------|----------|------------|--------|
| 0 | 0 | 22.0 | 7.2500 | 1 | False | False | True | False | |
| 1 | 1 | 38.0 | 71.2833 | 1 | True | False | False | True | |
| 2 | 1 | 26.0 | 7.9250 | 1 | False | False | True | False | |
| 3 | 1 | 35.0 | 53.1000 | 1 | True | False | False | False | |
| 4 | 0 | 35.0 | 8.0500 | 1 | False | False | True | False | |
| 4 | | | | | | | | | • |

In [25]:

```
test_df.isnull().sum()
```

Out[25]:

| PassengerId | 0 |
|--------------|-----|
| Pclass | 0 |
| Name | 0 |
| Sex | 0 |
| Age | 86 |
| SibSp | 0 |
| Parch | 0 |
| Ticket | 0 |
| Fare | 1 |
| Cabin | 327 |
| Embarked | 0 |
| dtype: int64 | |

In [26]:

```
test_data = test_df.copy()
test_data["Age"].fillna(train_df["Age"].median(skipna=True), inplace=True)
test_data["Fare"].fillna(train_df["Fare"].median(skipna=True), inplace=True)
test_data.drop('Cabin', axis=1, inplace=True)
test_data['TravelAlone']=np.where((test_data["SibSp"]+test_data["Parch"])>0,0,1)
test_data.drop('SibSp', axis=1, inplace=True)
test_data.drop('Parch', axis=1, inplace=True)
testing = pd.get_dummies(test_data, columns=["Pclass","Embarked","Sex"])
testing.drop('Sex_female', axis=1, inplace=True)
testing.drop('PassengerId', axis=1, inplace=True)
testing.drop('Name', axis=1, inplace=True)

final_test = testing
final_test.head()
```

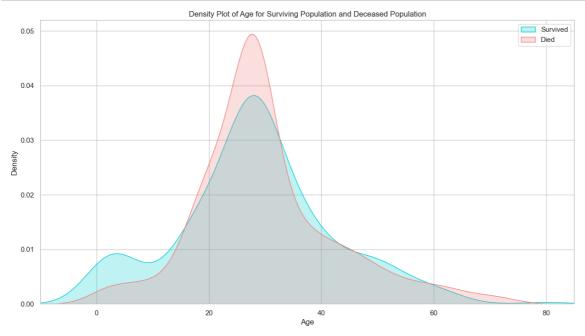
Out[26]:

| | Age | Fare | TravelAlone | Pclass_1 | Pclass_2 | Pclass_3 | Embarked_C | Embarked_Q | Em |
|---|------|---------|-------------|----------|----------|----------|------------|------------|----|
| 0 | 34.5 | 7.8292 | 1 | False | False | True | False | True | |
| 1 | 47.0 | 7.0000 | 0 | False | False | True | False | False | |
| 2 | 62.0 | 9.6875 | 1 | False | True | False | False | True | |
| 3 | 27.0 | 8.6625 | 1 | False | False | True | False | False | |
| 4 | 22.0 | 12.2875 | 0 | False | False | True | False | False | |
| 4 | | | | | | | | | • |

EXPLORATORY DATA ANALYSIS

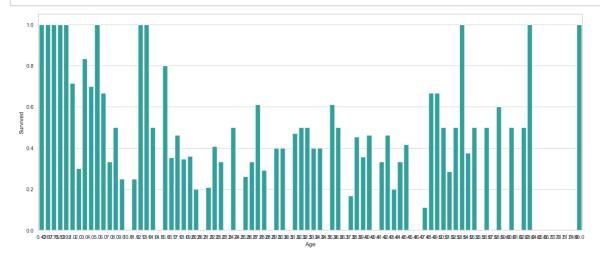
In [28]:

```
plt.figure(figsize=(15,8))
ax = sns.kdeplot(final_train["Age"][final_train.Survived == 1], color="darkturquoise", s
sns.kdeplot(final_train["Age"][final_train.Survived == 0], color="lightcoral", shade=Tru
plt.legend(['Survived', 'Died'])
plt.title('Density Plot of Age for Surviving Population and Deceased Population')
ax.set(xlabel='Age')
plt.xlim(-10,85)
plt.show()
```



In [29]:

```
plt.figure(figsize=(20,8))
avg_survival_byage = final_train[["Age", "Survived"]].groupby(['Age'], as_index=False).m
g = sns.barplot(x='Age', y='Survived', data=avg_survival_byage, color="LightSeaGreen")
plt.show()
```



```
In [30]:
```

```
final_train['IsMinor']=np.where(final_train['Age']<=16, 1, 0)</pre>
print(final_train['IsMinor'])
0
       0
1
       0
2
       0
3
       0
4
       0
886
       0
887
       0
       0
888
889
       0
890
       0
Name: IsMinor, Length: 891, dtype: int32
In [31]:
final_test['IsMinor']=np.where(final_test['Age']<=16, 1, 0)</pre>
print(final_test['IsMinor'])
       0
0
1
       0
2
       0
3
       0
4
       0
413
       0
414
       0
415
       0
       0
416
417
Name: IsMinor, Length: 418, dtype: int32
In [ ]:
```